Patterns Over Principles: The Fragility of Inductive Reasoning in LLMs under Noisy Observations

Anonymous ACL submission

Abstract

001 Inductive reasoning, a cornerstone of human cognition, enables generalization from limited data but hasn't yet been fully achieved by large language models (LLMs). While mod-005 ern LLMs excel at reasoning tasks, their ability to maintain stable and consistent rule abstraction under imperfect observations remains 007 underexplored. To fill this gap, in this work, we introduce Robust Rule Induction, a task that evaluates LLMs' capability in inferring rules from data that are fused with noisy examples. To address this task, we further propose Sample-steered Rule Refinement (SRR), a method enhancing reasoning stability via observation diversification and execution-guided feedback. Experiments across arithmetic, cryptography, and list functions reveal: (1) SRR 017 018 outperforms other methods with minimal performance degradation under noise; (2) Despite slight accuracy variation, LLMs exhibit instability under noise (e.g., 0% accuracy change with only 70% consistent score); (3) Counterfactual task gaps highlight LLMs' reliance on 024 memorized patterns over genuine abstraction. Our findings challenge LLMs' reasoning robustness, revealing susceptibility to hypothesis drift and pattern overfitting, while providing empirical evidence critical for developing human-like inductive systems.

1 Introduction

037

041

Inductive reasoning—the cognitive capacity to generalize from specific instances to universal principles, is fundamental to human intelligence (Lake et al., 2017). This ability enables humans to abstract latent patterns from sparse data while maintaining robustness against conflicting evidence (Feldman, 1997), as exemplified by children mastering linguistic rules despite exposure to occasional grammatical errors. The robustness of inductive reasoning manifests in resistance to pattern interference (preserving learned rules when encountering contradictory examples) and tolerance to observational imperfections (maintaining stable performance under noisy learning conditions). Understanding and quantifying this robustness becomes increasingly crucial as artificial systems approach complex real-world applications where clean data remains elusive. 042

043

044

047

048

053

054

056

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

076

078

079

081

Recent advances in large language models (LLMs) have demonstrated remarkable performance across various reasoning tasks (OpenAI, 2024b; DeepSeek-AI et al., 2024), reigniting interest in comparing machine and human reasoning paradigms (Collins et al., 2024). While contemporary studies showcase LLMs' proficiency in inductive reasoning, their capacity for robustness remains questionable. Unlike humans who can rapidly converge on correct rules through Bayesian hypothesis updating (Tenenbaum et al., 2011; Lake et al., 2015), current models often exhibit unstable reasoning trajectories when noise disrupts their thinking process (Zhou et al., 2024). A growing body of evidence suggests that LLMs primarily engage in pattern matching rather than genuine reasoning (Mirzadeh et al., 2024; Wu et al., 2024).

Despite increasing attention to LLM inductive reasoning capabilities, current works exhibit limitations that obscure true inductive reasoning robustness. First, in prevailing evaluations (Alet et al., 2021; Gendron et al., 2024), models just predict outputs for novel inputs given exemplars, bypassing explicit rule verification. It fails to diagnose where and why rule induction fails. Existing studies that do generate intermediate rules either disregard noisy learning conditions or assess performance solely through aggregate metrics like task accuracy gap. This neglects instance-level reasoning consistency-whether models maintain stable rule interpretations when exposed to conflicting patterns, which is a key indicator of human-like robust induction. Besides, prior methods often repurpose existing benchmarks (e.g., SCAN; Lake and Baroni

115

116

117

118

119

120

121

122

124

125

126

127

128

129

130

131

132

2018, ARC; Chollet 2019) without modification, where potential data contamination may undermine validity as models could use memorized solutions.

To address these gaps, we introduce a novel task of robust rule induction, which challenges models to identify underlying functions from input-output exemplars containing controlled noise injections. Our approach features three key innovations: (1) An evaluation pipeline from data synthesis to automatic evaluation, requiring models to explicitly output rules and enabling direct validation through programmatic execution. (2) A multi-dimensional robustness assessment with conventional accuracy and instance-level metrics. (3) A sample-steered refinement (SRR) methodology that improves both the capability and stability of LLMs' inductive reasoning through diverse sampling and diagnostic refinement.

We curate three datasets spanning arithmetic, cryptography, and list operations-domains requiring progressively abstract rule formalization. The proposed SRR substantially improves performance across most tasks and maintains stability compared to other methods. However, comprehensive experiments reveal a critical dichotomy: while the accuracy variance achieves less than 5.0%, the consistency metric exposes fundamental instability. Furthermore, LLMs exhibit significant performance inconsistency with the same task category, performing considerably better on more familiar tasks than others, despite comparable complexity. Correct predictions may stem from recited patterns rather than stable rule abstraction. These findings challenge the "real" reasoning capabilities of LLMs.

2 Preliminary

2.1 Problem Definitions

We formulate robust rule induction as the problem of identifying latent mapping functions from inputoutput pairs while maintaining consistent reasoning performance under noisy conditions. While prior works (Qiu et al., 2024; Wang et al., 2024) focus on basic rule discovery capabilities, we specifically investigate the robustness of inductive reasoning when the observations contain conflicting patterns. Some examples are shown in Figure 1. The observations may contain both normal and noisy examples, and the model is required to infer the correct rule regardless of the noise.

Formally, given a dataset $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$ where $x \in \mathcal{X}$ denotes inputs and $y \in \mathcal{Y}$ repre-



Figure 1: Example instances with noise and rules from Arithmetic_{base-7}, Cryptography_{Caesar} and List Functions.

sents outputs generated by an underlying function $f: \mathcal{X} \to \mathcal{Y}$, the objective is to induce an approximation \hat{f} from observed examples $\mathcal{D}_{\text{seen}} \subseteq \mathcal{D}$. Robust rule induction should satisfy: $\hat{f}(x) = f(x)$ for $(x, y) \in \mathcal{D}_{\text{unseen}}$ despite $\mathcal{D}_{\text{seen}}$ containing noisy samples $\mathcal{D}_{\text{noise}}$ where $y \neq f(x)$.

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

155

156

157

158

159

160

161

162

163

164

165

168

169

171

2.2 Evaluation Pipeline

As illustrated in figure 2, to evaluate the robustness of inductive reasoning in large language models, we propose a novel evaluation pipeline consisting of three main stages: data synthesis, model inference, and performance evaluation. In the data synthesis stage, due to the characteristics of our data, normal examples, noisy examples, and test examples for each instance in the dataset can be automatically generated by programs given the rule and noise definitions. The first two are mixed to form seen examples, while the latter is used to evaluate the model's reasoning ability. In the model inference stage, we prompt the language model with seen examples in standard input-output (IO) format, and the induced rules are restricted to Python function form for automatic evaluation. Finally, we evaluate the rule by executing it on the test examples. The instance is considered solved if all the test examples are correctly solved. More details about the evaluation pipeline, like the data synthesis process, can be found in appendix A.

We use task accuracy, which is the proportion of solved instances over the total number of instances, and its change under different conditions as the evaluation metrics. While task accuracy is a comprehensive metric to evaluate the model's reasoning ability, it may not fully capture the robustness of the model at the instance level. To address this, we introduce the consistency score, which is defined as follows:

Consistency Score :=
$$\frac{1}{N} \sum_{i=1}^{N} \mathbb{I}[Sol_i^c = Sol_i^n]$$
 170

where N is the total number of instances in the



Stage 3: Performance Evaluation

Stage 2: Model Inference

Figure 2: Evaluation pipeline exemplified by base-9 addition, consisting of three stages: (1) Data Synthesis, generating normal, noisy and test examples; (2) Model Inference, prompting models with seen examples to induce rules in Python function form; (3) Performance Evaluation, executing induced rules on test examples to assess correctness and robustness under noise.

dataset, $\mathbb{I}[\cdot]$ is the indicator function, Sol_i^c is whether the *i*-th instance is solved without noise, 1 if solved and 0 otherwise, and Sol_i^n is whether the *i*-th instance is solved with noise. This metric quantifies the model's ability to maintain stable reasoning conclusions under noise and offers a more granular view of the model's robustness.

172

173

174

178

179

180

181

188

3 Sample-steered Rule Refinement

To address noisy rule induction, we propose Sample-steered Rule Refinement (SRR), a novel framework that combines hypothesis generation, contrastive evidence sampling, and iterative selfcorrection. This method conducts inductive reasoning through three-phase optimization: (1) contrastive hypothesis generation to bootstrap diverse rule candidates, (2) diagnostic feedback construction through evidence-aware sampling, and (3) iterative rule refinement using failure-driven corrections. The full algorithm is detailed in Algorithm 1.

Diversity-aware Hypothesis Generation We 191 first generate K hypotheses from random subsets 192 of the seen examples plus one hypothesis using the full observations. This step ensures coverage 194 of both local patterns and global consistency. The number of examples in the subset is less than the full seen set, thus the noise in the subset has a 197 greater impact. The sampled subset may also con-198 tain completely noise-free examples, which can 199 better distinguish the rules.

Diagnostic Feedback Construction After generating the initial hypotheses, the induced rule is exe-202

Algorithm 1 SRR Framework

```
Require: Seen examples \mathcal{D}_{\text{seen}}, LLM M, max iterations T,
    max subsets K, threshold \tau
Ensure: Best hypothesis \hat{f}^* (in Python function form)
     \mathcal{H}_0 \leftarrow \emptyset
    for k = 1 to K do
           \mathcal{D}_k \leftarrow \text{SampleSubset}(\mathcal{D}_{\text{seen}})
           \mathcal{H}_0 \leftarrow \mathcal{H}_0 \cup \{ M(\mathcal{D}_k, "Generate rule function") \} 
     end for
     \mathcal{H}_0 \leftarrow \mathcal{H}_0 \cup \{M(\mathcal{D}_{\text{seen}}, \text{"Generate rule function"})\}
     \hat{f}_0 \leftarrow \arg \max_{h \in \mathcal{H}_0} \operatorname{Acc}(h, \mathcal{D}_{\operatorname{seen}})
    for t = 1 to T do
           if \operatorname{Acc}(\hat{f}_{t-1}, \mathcal{D}_{\operatorname{seen}}) \geq \tau then
                return \hat{f}_{t-1}
           end if
           C_t \leftarrow \text{CorrectExamples}(\hat{f}_{t-1}, \mathcal{D}_{\text{seen}})
           \mathcal{E}_t \leftarrow \text{WrongExamples}(\hat{f}_{t-1}, \mathcal{D}_{\text{seen}})
           Feedback \leftarrow \{C_t[:n], \mathcal{E}_t[:m]\} /* \text{ Sample n,m cases }*/
           h_{\text{new}} \leftarrow M(\hat{f}_{t-1}, \text{Feedback}, "\text{Revise rule"})
           if \operatorname{Acc}(h_{\operatorname{new}}, \mathcal{D}_{\operatorname{seen}}) > \operatorname{Acc}(\hat{f}_{t-1}, \mathcal{D}_{\operatorname{seen}}) then
                 f_t \leftarrow h_{\text{new}}
           else
                 f_t \leftarrow f_{t-1}
           end if
    end for
    return \hat{f}^* = \arg \max_t \operatorname{Acc}(\hat{f}_t, \mathcal{D}_{\text{seen}})
```

cuted on the seen examples via an external Python interpreter. We then collect the correct and incorrect cases and sample these cases as feedback. We focus more on the incorrect cases, as they provide more information about the rule refinement, while the correct cases are used as positive feedback to reinforce the rule. The hypothesis with the highest accuracy on the seen examples is selected as the initial selected rule.

Iterative Rule Refinement At each iteration, the model receives the current selected rule with for-

212

213

203

3

matted feedback to generate a new rule. The new hypothesis is compared to the previous one, and the more accurate one is selected. This iterative process continues until the accuracy of the rule on the seen examples exceeds a predefined threshold or the maximum number of iterations is reached.

4 Experiments

220

222

228

234

236

240

241

242

243

245

247

249

In this section, we evaluate the robustness of inductive reasoning in language models on different tasks. We also compare the performance of the models under different noise levels and analyze the effectiveness of the different methods in enhancing the robustness of inductive reasoning.

4.1 Experimental Setup

Consistent with previous studies (Wang et al., 2024; Qiu et al., 2024), we adopt few-shot prompting to assess the models' inductive capabilities. Each instance contains 10 normal examples, 5 noisy examples, and 10 test examples. The normal and noisy examples are mixed to form 10 seen examples, which serve as prompts. This approach preserves task semantics while introducing controlled perturbations, thereby simulating real-world scenarios where observational data often contains inherent imperfections. LLMs are explicitly informed that examples may contain some noise. We formulate the output of LLMs as Python functions and execute them on the test examples to automatically evaluate the inferred rules. More details about the experiments can be found in the appendix B.

4.2 Datasets

We use three datasets with their subsets: Arithmetic, Cryptography, and List Functions. These datasets offer different rule induction challenges, including mathematical calculations, symbol representation, and list operations. Comprehensive statistics are shown in Table 1.

Dataset	Subset	# Tasks	# Normal	# Noise	# Test
Arith.	7-base 8-base 9-base	$ \begin{array}{c} 100 \\ 100 \\ 100 \end{array} $	1000 1000 1000	500 500 500	$1000 \\ 1000 \\ 1000$
Crypto.	Caesar Atbash Keyboard	100 100 100	1000 1000 1000	500 500 500	1000 1000 1000
List Func.	N/A	250	2500	1250	2500

Table 1: The number of tasks and examples per dataset. Arith., Crypto., List Func. denote Arithmetic, Cryptography and List Functions respectively.

Arithmetic The arithmetic task, proposed by (Wu et al., 2024), is a counterfactual two-digit addition task. Instead of the common base-10 system, the task uses base-8, 9, and 11 as the counterfactual setup to control the difficulty and avoid the impact of memorization over reasoning. In this task, we focus on the base-7, 8, 9 systems, and the noisy example is the common 10-based equations. The input of each example is two two-digit numbers in the corresponding base, connected by a plus sign, and the output is the sum of the two numbers. To ensure the model can reason correctly, we guarantee that there must be a carry-over in the addition process so that the model must reason about the carry-over instead of simply adding the numbers. 251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

266

267

268

269

270

271

272

273

274

275

276

277

278

279

281

283

284

285

286

287

289

291

292

293

294

295

296

297

298

299

300

Cryptography We use three types of substitution ciphers: Caesar, Atbash, and Keyboard. The Caesar cipher is a cipher that shifts the alphabet by a fixed number of positions. The Atbash cipher is a cipher that replaces each letter with the letter symmetrically opposite in the alphabet (e.g. $A \rightarrow Z$). The Keyboard cipher replaces letters according to their positions in the alphabet with the corresponding positions on the keyboard (e.g. $A \rightarrow Q$, $B \rightarrow W$). The input of each example is a word, and the output is the encrypted text. We randomly replace the characters in the output with other characters to generate noisy examples.

List Functions The list functions dataset (Rule, 2020) evaluates the concept learning ability in the domain of cognitive science. The task is to induce a function that maps a list of numbers to another list of numbers. The dataset has 250 tasks. Each item in the dataset corresponds to a list manipulation operation, such as sorting, reversing, or filtering. We randomly generate the examples of each task under controlled conditions to ensure the rule can be induced. The noisy examples are generated by randomly replacing the numbers in the output list with other numbers.

4.3 Robustness Under Different Noise Levels

We first evaluate the robustness of inductive reasoning in language models under different noise injection ratios. Noise levels are defined as the proportion of noisy examples in the seen examples. We test three representative models: GPT-40-mini (OpenAI, 2024a), GPT-40 (OpenAI, 2024b), and DeepSeek-V3 (DeepSeek-AI et al., 2024) and ask the models to directly infer the rules from the seen examples without any additional output. Since



Figure 3: Consistency score(%) with clean data of different models on the Cryptography and List Functions datasets under different noise levels.

Model	(List Functions		
	Caesar	Atbash	Keyboard	
GPT-40-mini	$28.6_{\pm 1.2}$	$6.3_{\pm 1.2}$	$0.0_{\pm 0.0}$	$31.3_{\pm 0.2}$
10% noise	$32.0_{\pm 2.4}$	$12.7_{\pm 2.5}$	$3.7_{\pm 0.5}$	28.4 ± 1.4
20% noise	$29.0_{\pm 0.8}$	13.3 ± 2.6	$1.3_{\pm 1.2}$	23.4 ± 1.0
30% noise	21.3 ± 0.9	$16.0_{\pm 0.0}$	$3.3_{\pm 0.5}$	$19.7_{\pm 0.2}$
GPT-40	68.3 ± 0.5	48.0 ± 2.2	$4.7_{\pm 0.5}$	$39.8_{\pm 0.5}$
10% noise	$75.3_{\pm 0.5}$	$34.6_{\pm 0.9}$	$3.7_{\pm 0.5}$	$36.9_{\pm 1.0}$
20% noise	$75.3_{\pm 2.1}$	27.3 ± 0.5	$1.7_{\pm 0.9}$	$31.9_{\pm 0.8}$
30% noise	$64.7_{\pm 4.7}$	$22.7_{\pm 2.4}$	$2.3_{\pm 1.2}$	$26.1_{\pm 0.5}$
DeepSeek-V3	$50.0_{\pm 1.6}$	$21.7_{\pm 2.5}$	25.0 ± 0.8	$43.1_{\pm 0.4}$
10% noise	$24.3_{\pm 1.7}$	15.3 ± 2.4	$15.7_{\pm 1.2}$	$40.8_{\pm 0.9}$
20% noise	$15.3_{\pm 3.7}$	$21.3_{\pm 2.5}$	$11.3_{\pm 1.2}$	$35.3_{\pm 1.9}$
30% noise	$11.7_{\pm 0.8}$	$16.7_{\pm 0.5}$	$11.3_{\pm 3.1}$	$30.5_{\pm 0.7}$

Table 2: Task accuracy (%) on Cryptography (including Caesar, Atbash, and Keyboard subtasks) and List Functions under different noise levels (proportion of noise in seen examples). Results are shown in mean \pm standard deviation over 3 independent runs.

the models cannot solve the Arithmetic task under the Direct Output (DO) setting, we present the results for the Cryptography and List Functions. The results are shown in Table 2. To better evaluate the robustness of the models, we further investigate the consistency score with clean data under different noise levels, as shown in Figure 3. According to the results, we have the following observations.

301

302

306

307

310

312

313

315

317

319

323

First, contrary to conventional assumptions, noise introduction does not universally degrade performance, instead, the models exhibit performance fluctuations, demonstrating their inherent sensitivity to conflicting patterns. For example, GPT-40 achieves improved accuracy on Caesar cipher tasks at 10% noise (7.0% absolute improvement over clean data), while List Functions exhibit monotonic performance decay with increasing noise levels. In some cases, moderate noise improves performance, consistent with the findings of Zhou et al. (2024).

With the exception of tasks where models fundamentally struggle, experimental results demonstrate **a decline in consistency scores as noise lev-** els escalate, and it declines more sharply than its task accuracy variation. This discrepancy indicates noise introduces bidirectional reasoning instability: models not only fail on previously solvable instances (noise interference) but also succeed on originally challenging cases (incorrect generalization), leading to different performance changes. 324

325

326

327

328

329

330

332

333

334

335

336

339

340

341

344

345

346

347

348

349

350

352

354

356

357

360

4.4 Method-wise Effectiveness Comparison

Direct Output may limit the model's reasoning ability, as the model must output the rule directly without any intermediate steps. To systematically assess the robustness of different reasoning paradigms, we compare our Sample-steered Rule Refinement (SRR) method with three reasoning methods: (1) Chain of Thought (CoT; Wei et al. 2022), which decomposes reasoning into step-bystep rationales; (2) Self-Consistency (SC; Wang et al. 2023), which aggregates multi CoT trajectories through majority voting; and (3) Self-Refine (SR; Madaan et al. 2023), which iteratively improves hypotheses using self-generated feedback. Table 3 presents the task accuracy under 10% noise and the deviation from clean data for GPT-40 and DeepSeek-V3.

Superior Performance of SRR As shown in Table 3, our SRR framework achieves state-ofthe-art performance across 13/14 task-model combinations while exhibiting minimal performance degradation (2.1% average drop vs. 3.6%-8.5% for baselines). This dual advantage stems from two mechanisms: (1) *Diversity-aware hypothesis generation* explores broader solution spaces through subset sampling, outperforming SC's majority voting that amplifies the similar patterns (Table 4); (2) *Execution-guided feedback* leverages Python interpreters for objective error detection, circumventing LLMs' inherent deductive limitations (Chen et al.,

Model	Method	Arithmetic		Cryptography			List Functions	
		7-base	8-base	9-base	Caesar	Atbash	Keyboard	
CDT 4a	DO CoT	$\begin{array}{c} 0.0\\ 3.0(\leftrightarrow \underline{0.0})\\ 0.0(\pm 1.0)\end{array}$	0.0 $8.0(\downarrow 14.0)$	0.0 $8.0(\uparrow 5.0)$	$\begin{array}{c c} 75.3(\uparrow 7.0) \\ 85.0(\uparrow 1.0) \\ 85.0(\downarrow 0.0) \end{array}$	$34.6(\downarrow 13.4)$ $20.0(\downarrow 7.0)$	$3.7(\downarrow 1.0) \\ 4.0(\leftrightarrow \underline{0.0}) \\ 5.0(\downarrow 2.0)$	$\begin{vmatrix} 36.9(\downarrow 2.9) \\ 40.4(\downarrow 5.2) \\ 42.4(\downarrow 2.2) \end{vmatrix}$
GP1-40	SC SR Ours	$0.0(\downarrow 1.0)$ $1.0(\downarrow 4.0)$ $5.0(\downarrow 1.0)$	$18.0(\downarrow 1.0)$ $18.0(\downarrow 4.0)$ $51.0(\uparrow 6.0)$	$4.0(\downarrow 3.0)$ $4.0(\downarrow 2.0)$ $19.0(\downarrow 2.0)$	85.0 ($\leftrightarrow \underline{0.0}$) 81.0($\downarrow 2.0$) 85.0 ($\uparrow 3.0$)	$29.0(\downarrow 11.0) 23.0(\downarrow 8.0) 52.0(\downarrow 3.0)$	$5.0(\downarrow 2.0)$ $5.0(\leftrightarrow 0.0)$ $9.0(\uparrow 1.0)$	$\begin{vmatrix} 42.4(\downarrow 3.2) \\ 39.6(\downarrow 4.4) \\ 57.2(\downarrow 1.6) \end{vmatrix}$
Deepseek-V3	DO CoT SC SR Ours	$\begin{array}{c} 0.0 \\ 77.5(\downarrow 6.0) \\ 83.0(\downarrow 3.0) \\ 70.0(\downarrow 10.0) \\ \textbf{96.0}(\downarrow \underline{1.0}) \end{array}$	$\begin{array}{c} 0.0\\ 83.0(\downarrow\ 13.0)\\ 93.0(\downarrow\ 6.0)\\ 74.0(\downarrow\ 9.0)\\ \textbf{95.0}(\downarrow\ \underline{4.0})\end{array}$	$\begin{array}{c} 0.0 \\ 67.5(\downarrow 14.0) \\ 81.0(\downarrow 3.0) \\ 68.0(\uparrow 6.0) \\ \textbf{94.0}(\leftrightarrow \underline{0.0}) \end{array}$	$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	$\begin{array}{c} 15.3(\downarrow 6.4)\\ 26.0(\downarrow 5.5)\\ 40.0(\downarrow 3.0)\\ 19.0(\downarrow \underline{1.0})\\ \textbf{52.0}(\downarrow \underline{1.0}) \end{array}$	$\begin{array}{c} \textbf{15.7}(\downarrow 9.3) \\ 5.0(\downarrow 8.5) \\ 7.0(\downarrow 4.0) \\ 8.0(\downarrow 3.0) \\ 11.0(\downarrow \underline{2.0}) \end{array}$	$ \begin{vmatrix} 40.8(\downarrow \underline{2.3}) \\ 52.0(\downarrow 6.4) \\ 56.0(\downarrow 3.2) \\ 47.2(\downarrow 10.0) \\ 64.8(\downarrow 2.8) \end{vmatrix} $

Table 3: Task accuracy (%) on different datasets under 10% noise. The numbers in parentheses are the change compared to the clean data, and the arrows indicate the direction of the change. **Bold** indicates the best performance, and <u>underline</u> indicates the smallest change.

Dataset	GI	РТ-40	DeepSeek-V3		
	SC	SRR-0	SC	SRR-0	
Arithmetic ₇	0.0	4.0	83.0	95.0	
Arithmetic ₈	6.0	40.0	93.0	95.0	
Arithmetic ₉	0.0	15.0	94.0	91.0	
Crypto _{Caesar}	85.0	85.0	86.0	86.0	
Crypto _{Atbash}	29.0	51.0	40.0	48.0	
Crypto _{Keyboard}	5.0	8.0	7.0	9.0	
List Functions	42.4	54.0	56.0	62.8	

Table 4: Task accuracy (%) under 10% noise of SC and the initial rule in SRR (SRR-0) on different datasets.

2023b; Cheng et al., 2024).

361

365

367

370

371

372

374

377

385

Consistency Scores Reveal Hidden Instability While task accuracy provides comprehensive insights, consistency scores uncover fundamental reasoning fragility. As depicted in Figure 4, the Atbash cipher task exhibits particularly low consistency despite modest accuracy changes. The slight variation in accuracy may merely be an illusion created by the combined effects of noise interference (solved→unsolved cases) and incorrect generalization (unsolved→solved cases). Prior and concurrent work's singular focus on accuracy fluctuations (Qiu et al., 2024; Zhou et al., 2024; Mirzadeh et al., 2024; Huang et al., 2025a) overlooks this critical duality in reasoning robustness.

Counterfactual Challenges Expose Knowledge Reliance We observe dramatic performance gaps on counterfactual tasks like base-7 and base-9 arithmetic (90%+ vs. 10%- for DeepSeek-V3 vs. GPT-40) and Keyboard ciphers. This observation is consistent with the findings reported in Wu et al. (2024). Diagnostic analysis of responses from LLMs shows models default to familiar templates rather than true induction. GPT-40 persistently misinterprets base-7 and base-9 addition as "base-



Figure 4: Consistency score (%) between clean data and data with 10% noise of DeepSeek-V3.

387

388

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

8" or "decimal sum with constant", while Atbash and Keyboard ciphers get erroneously classified as Caesar shifts in thinking process or self-generated feedback. Although Atbash and Caesar cipher exhibit similar levels of complexity in transformation, their performance differences are significant. These failures reveal that there is a pattern overfitting in the reasoning process. Current models' inductive reasoning essentially operates through pattern matching rather than abstract induction. When the scarcity of counterfactual tasks in pretraining data forces models to rely on genuine rule induction rather than recitation, model performance plummets dramatically.

4.5 Extended Explorations

To further investigate the robustness of LLMs' inductive reasoning, we conduct a comparative evaluation against DeepSeek-R1 (DeepSeek-AI et al., 2025), one of the state-of-the-art reasoning models, and human reasoning patterns.

Comparion with DeepSeek-R1 We focus on the 406 Atbash cipher and List Functions. Atbash pro-407 vides unified task semantics yet challenges models 408 with unfamiliar transformation logic, while List 409 Functions captures diverse rule abstraction scenar-410 ios. Table 5 compares task accuracy and consis-411 tency scores under clean and 10% noise conditions 412 with DeepSeek-V3 (SRR). Notably, DeepSeek-R1 413 achieves higher task accuracy on both tasks. How-414 ever, its consistency scores are not competitive 415 enough, indicating unresolved instability. Detailed 416 breakdowns of consistency scores are shown in 6. 417 For Atbash, the comparable number of right-to-418 wrong and wrong-to-right suggests instability and 419 randomness in its reasoning process. Manual in-420 spection shows DeepSeek-R1 still interprets At-421 bash cipher as character shifts in failed cases, mir-422 roring previous pattern-overfitting behavior. For 423 List Functions, the elevated right-to-wrong rate in-424 dicates its sensitivity to input perturbations.

Dataset	Model	Clean Acc	Noise Acc	Consistency
Atbash	SRR DeepSeek-R1	$53.0 \\ 65.0$	$\begin{array}{c} 52.0(-1.0) \\ 65.0(\pm 0.0) \end{array}$	69.0 70.0
List Func.	SRR DeepSeek-R1	67.6 76.0	$\begin{array}{c} 64.8(-2.8) \\ 68.4(-7.6) \end{array}$	90.8 87.6

Table 5: Task accuracy (%) and consistency score (%) between DeepSeek-V3 with SRR and DeepSeek-R1 under clean and noisy conditions.

Dataset	Model	# BothR	# BothW	# RtoW	# WtoR
Atbash	SRR	37	32	16	15
	DeepSeek-R1	50	20	15	15
List Func.	SRR	154	73	15	8
	DeepSeek-R1	165	54	25	6

Table 6: Detailed breakdowns of consistency score. BothR, BothW, RtoW, WtoR represent both right, both wrong, right to wrong from clean condition to noisy condition and wrong to right respectively.

425

426

427

428

429

430

431

432 433

434

435

436

437

Human Reasoning Comparison Rule (2020) reports human performance on List Functions. We divide the tasks into three difficulty levels in a 5:3:2 ratio based on the sorted mean human performance. To contrast with human-like thinking patterns, we analyze consistency across 12 trials (4 noise levels \times 3 runs) in Section 4.3. Figure 5 visualizes the distribution of consistency and performance per task in List Functions, stratified by difficulty. LLMs show macro-level alignment yet micro-level divergence compared to human reasoning. While LLMs broadly mirror human-like stability, showing higher consistency on both simple and hard tasks, with moderate instability on medium-difficulty tasks, their internal patterns reveal critical deviations. For hard tasks, models display unpredictability (e.g., sporadic success) rather than systematic incapacity. Simple tasks, despite their low complexity, exhibit in consistent performance unrelated to intrinsic difficulty. Notably, GPT-40 and DeepSeek-V3 demonstrate similar behavioral trends, suggesting a shared inductive bias with humans. 438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467



Figure 5: Task-solving consistency of DeepSeek-V3 and GPT-40 on List Functions. Each cell represents a task, arranged by ascending difficulty (top-to-bottom, left-to-right). Colors denote correctness patterns: 12R (all correct) to 12W (all wrong), with intermediate states (e.g., 1W11R: 1 wrong, 11 correct).

5 Discussion

Our experimental results reveal limitations in LLMs' inductive reasoning. We try to analyze the results in this section.

Process Analysis Inductive reasoning can be conceptualized through a Bayesian paradigm (Tenenbaum et al., 2011), where models update posterior distributions over a hypothesis space based on observations. The introduction of noise disrupts the reasoning through dual mechanisms: **Noise-Induced Hypothesis Drift**, where the initial hypothesis space becomes misaligned with true rules when noise introduces conflicting patterns. This drift particularly impacts Direct Output (DO) method, which lacks intermediate reasoning steps and relies heavily on the initial hypothesis space. This is more evident for tasks with similar internal patterns, the model's performance fluctuates significantly (Cryptography) or even fails (Arithmetic)

compared to other methods. Evidence Ambiguity 468 Amplification, where noise reduces the effective 469 signal-to-noise ratio during posterior optimization. 470 Methods with iterative or step-by-step reasoning 471 suffer error accumulation cascades-each reason-472 ing step propagates uncertainty. As shown in Ta-473 ble 3, CoT and SR exhibit larger average accuracy 474 drops and worse performance. SRR avoids it by 475 objectively identifying errors through execution. 476

Implications for Robust Inductive Reasoning 477 LLMs heavily rely on prior knowledge rather than 478 pure induction. They blend memorized pattern 479 matching with shallow reasoning. When observa-480 tions align with priors (e.g., Caesar cipher), they 481 482 demonstrate pseudo-robustness. However, conflicting patterns (noise) and counterfactual scenarios 483 (unseen rules) expose this fragility-models either 484 default to familiar templates or enter unstable hy-485 pothesis oscillations. Even the state-of-the-art rea-486 soning model exhibits unstable performance in uni-487 fied tasks. This contrasts sharply with human few-488 shot learning, where true rule abstraction enables 489 stable generalization (Lake et al., 2015). Achieving 490 human-level robust reasoning and few-shot induc-491 tion requires disentangling knowledge recitation 492 from inductive rule formation. 493

6 Related Work

494

495

496

497

498

500

501

502

503

504

505

506

507

509

510

511

512

513

514

515

517

Inductive Reasoning The study of inductive reasoning has been a long-standing focus across multiple disciplines. Early work (Heit, 2000) established foundational properties of inductive reasoning. In cognitive science, induction is considered a process of probabilistic belief updating within the Bayesian paradigm (Tenenbaum et al., 2011). Human cognition and learning are explained through the integration of prior knowledge with observed data to compute posterior distributions. Comparative studies (Lake et al., 2015, 2017) have highlighted the contrast between human learners and machine intelligence. With the advent of pre-trained language models (Brown et al., 2020), research on inductive reasoning has shifted from domain-specific and neural formulation (Tian et al., 2020; Odena et al., 2021; Sablé-Meyer et al., 2022) to natural language. Initial approaches (Alet et al., 2021; Gendron et al., 2024; Mirchandani et al., 2023; Yang et al., 2024) predominantly relied on input-output (IO) prompting, which evaluates model performance on unseen examples without explicit rule articulation. However, this paradigm overlooks the internal ruleinference process, as it conflates rule induction with rule execution capabilities. Recent efforts (Wang et al., 2024) align more closely with our work by explicitly generating intermediate rules. Qiu et al. (2024) proposed a thorough evaluation containing noisy conditions, whose findings partially overlap with ours, but their analysis is confined to a single dataset and prioritized accuracy metrics over consistency. We further explore the challenges posed by counterfactual tasks and uncover the limitations of LLMs in inductive reasoning.

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

Robustness of Reasoning in LLMs In this work, robustness refers to the ability to maintain consistent performance under imperfections or counterfactual scenarios (Elazar et al., 2021), which is intrinsically linked to the diversity and unpredictability of generation process (Zhang et al., 2023; Huang et al., 2025b). Prior works into LLM behavioral variability exhibit inconsistent performance across time (Tu et al., 2024; Chen et al., 2023a). Recent works focus on understanding their internal reasoning mechanisms by altering conditions, such as changing numerical values in mathematical tasks (Mirzadeh et al., 2024; Huang et al., 2025a) or examining performance on counterfactual tasks (Wu et al., 2024). Zhou et al. (2024) investigates the impact of noisy rationales on model performance. These studies primarily measure robustness through global accuracy changes, overlooking inconsistencies at the instance level. Our work addresses this gap by introducing a consistency score that quantifies intra-task stability, providing a more granular view of model robustness.

7 Conclusion

In this paper, we explore the robustness of inductive reasoning in large language models under imperfect observations. Through introducing the Robust Rule Induction task with metrics of both holistic and individual levels, we systematically evaluate the ability of LLMs to maintain stable and consistent rule abstraction. The Sample-steered Rule Refinement outperforms other reasoning paradigms by effectively leveraging diversity-aware hypothesis generation and execution-guided feedback. Our findings reveal that while LLMs can demonstrate impressive reasoning capabilities, they are inherently sensitive to noise and prone to hypothesis drifting and pattern overfitting.

669

670

671

672

673

674

675

676

677

Limitations

566

586

588

591

593

594

595

596

597

602

604

605

606

607

610

611 612

613

614

615

616

We discuss the limitations of our work here: Our evaluation focuses on highly formalized and sym-568 bolic tasks. Real-world inductive reasoning often 569 involves ambiguous rules like social norms from text or visual patterns, which are not captured by 571 our formalism. The tasks in Arithmetic and Cryptography may be relatively simple and lack diversity. Expanding task diversity could reveal deeper 574 limitations in LLMs' capabilities. A large-scale assessment of human performance across all tasks 576 under varying levels of noise is beyond the scope of 577 this study, which poses certain limitations in comparing human reasoning patterns. The purpose of our study is to explore the robustness of LLMs' in-580 ductive reasoning capabilities, so we do not specif-581 ically tune the hyperparameters and prompt tem-582 plates. 583

Ethics Statement

The datasets we used are all publicly available, and our research does not involve any personal information. All data is generated by programs. Therefore, we anticipate that this paper does not raise any ethical concerns. We use ChatGPT to paraphrase some sentences.

References

- Ferran Alet, Javier Lopez-Contreras, James Koppel, Maxwell Nye, Armando Solar-Lezama, Tomas Lozano-Perez, Leslie Kaelbling, and Joshua Tenenbaum. 2021. A large-scale benchmark for few-shot program induction and synthesis. In Proceedings of the 38th International Conference on Machine Learning, volume 139 of Proceedings of Machine Learning Research, pages 175–186. PMLR.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In Proceedings of the 34th International Conference on Neural Information Processing Systems, NIPS '20, Red Hook, NY, USA. Curran Associates Inc.
- Lingjiao Chen, Matei Zaharia, and James Zou. 2023a. How is chatgpt's behavior changing over time? Preprint, arXiv:2307.09009.

- Xinyun Chen, Maxwell Lin, Nathanael Schärli, and Denny Zhou. 2023b. Teaching large language models to self-debug. Preprint, arXiv:2304.05128.
- Kewei Cheng, Jingfeng Yang, Haoming Jiang, Zhengyang Wang, Binxuan Huang, Ruirui Li, Shiyang Li, Zheng Li, Yifan Gao, Xian Li, Bing Yin, and Yizhou Sun. 2024. Inductive or deductive? rethinking the fundamental reasoning abilities of llms. Preprint, arXiv:2408.00114.
- François Chollet. 2019. On the measure of intelligence. Preprint, arXiv:1911.01547.
- Katherine M. Collins, Ilia Sucholutsky, Umang Bhatt, Kartik Chandra, Lionel Wong, Mina Lee, Cedegao E. Zhang, Tan Zhi-Xuan, Mark Ho, Vikash Mansinghka, Adrian Weller, Joshua B. Tenenbaum, and Thomas L. Griffiths. 2024. Building machines that learn and think with people. Nature Human Behaviour, 8(10):1851-1863. Publisher: Nature Publishing Group.
- DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Meng Li, Miaojun Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, R. J. Chen, R. L. Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan, S. S. Li, Shuang Zhou, Shaoqing Wu, Shengfeng Ye, Tao Yun, Tian Pei, Tianyu Sun, T. Wang, Wangding Zeng, Wanjia Zhao, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, W. L. Xiao, Wei An, Xiaodong Liu, Xiaohan Wang, Xiaokang Chen, Xiaotao Nie, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, X. Q. Li, Xiangyue Jin, Xiaojin Shen, Xiaosha Chen, Xiaowen Sun, Xiaoxiang Wang, Xinnan Song, Xinyi Zhou, Xianzu Wang, Xinxia Shan, Y. K. Li, Y. Q. Wang, Y. X. Wei, Yang Zhang, Yanhong Xu, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Wang, Yi Yu, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yuan Ou, Yuduan Wang, Yue Gong, Yuheng Zou, Yujia He, Yunfan Xiong, Yuxiang Luo, Yuxiang You,

Yuxuan Liu, Yuyang Zhou, Y. X. Zhu, Yanhong Xu, Yanping Huang, Yaohui Li, Yi Zheng, Yuchen Zhu, Yunxian Ma, Ying Tang, Yukun Zha, Yuting Yan, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Zizheng Pan, Zhen Huang, Zhipeng Xu, Zhongyu Zhang, and Zhen Zhang. 2025. Deepseek-r1: Incentivizing reasoning capability in Ilms via reinforcement learning. *Preprint*, arXiv:2501.12948.

678

679

681

688

700

701

702

703

704

705

709

710

711

712

713

714

715

716

717

718

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

740

- DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Haowei Zhang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Li, Hui Qu, J. L. Cai, Jian Liang, Jianzhong Guo, Jiaqi Ni, Jiashi Li, Jiawei Wang, Jin Chen, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, Junxiao Song, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Lei Xu, Leyi Xia, Liang Zhao, Litong Wang, Liyue Zhang, Meng Li, Miaojun Wang, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Mingming Li, Ning Tian, Panpan Huang, Peiyi Wang, Peng Zhang, Qiancheng Wang, Qihao Zhu, Qinyu Chen, Qiushi Du, R. J. Chen, R. L. Jin, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, Runxin Xu, Ruoyu Zhang, Ruyi Chen, S. S. Li, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shaoqing Wu, Shengfeng Ye, Shengfeng Ye, Shirong Ma, Shiyu Wang, Shuang Zhou, Shuiping Yu, Shunfeng Zhou, Shuting Pan, T. Wang, Tao Yun, Tian Pei, Tianyu Sun, W. L. Xiao, Wangding Zeng, Wanjia Zhao, Wei An, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, X. Q. Li, Xiangyue Jin, Xianzu Wang, Xiao Bi, Xiaodong Liu, Xiaohan Wang, Xiaojin Shen, Xiaokang Chen, Xiaokang Zhang, Xiaosha Chen, Xiaotao Nie, Xiaowen Sun, Xiaoxiang Wang, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xingkai Yu, Xinnan Song, Xinxia Shan, Xinyi Zhou, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, Y. K. Li, Y. Q. Wang, Y. X. Wei, Y. X. Zhu, Yang Zhang, Yanhong Xu, Yanhong Xu, Yanping Huang, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Li, Yaohui Wang, Yi Yu, Yi Zheng, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Ying Tang, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yu Wu, Yuan Ou, Yuchen Zhu, Yuduan Wang, Yue Gong, Yuheng Zou, Yujia He, Yukun Zha, Yunfan Xiong, Yunxian Ma, Yuting Yan, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Z. F. Wu, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhen Huang, Zhen Zhang, Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhibin Gou, Zhicheng Ma, Zhigang Yan, Zhihong Shao, Zhipeng Xu, Zhiyu Wu, Zhongyu Zhang, Zhuoshu Li, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Ziyi Gao, and Zizheng Pan. 2024. Deepseek-v3 technical report. Preprint, arXiv:2412.19437.
- Yanai Elazar, Nora Kassner, Shauli Ravfogel, Abhilasha Ravichander, Eduard Hovy, Hinrich Schütze, and Yoav Goldberg. 2021. Measuring and improving consistency in pretrained language models. *Transactions of the Association for Computational Linguistics*, 9:1012–1031.

741

742

743

744

745

747

749

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

765

766

767

768

769

770

771

772

773

774

775

776

777

778

779

780

781

782

783

784

785

786

787

788

789

790

791

792

793

794

795

- Jacob Feldman. 1997. The structure of perceptual categories. *Journal of mathematical psychology*, 41 2:145–70.
- Gaël Gendron, Qiming Bao, Michael Witbrock, and Gillian Dobbie. 2024. Large language models are not strong abstract reasoners. In *Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence, IJCAI-24*, pages 6270–6278. International Joint Conferences on Artificial Intelligence Organization. Main Track.
- Evan Heit. 2000. Properties of inductive reasoning. *Psychonomic Bulletin & Review*, 7(4):569–592.
- Kaixuan Huang, Jiacheng Guo, Zihao Li, Xiang Ji, Jiawei Ge, Wenzhe Li, Yingqing Guo, Tianle Cai, Hui Yuan, Runzhe Wang, Yue Wu, Ming Yin, Shange Tang, Yangsibo Huang, Chi Jin, Xinyun Chen, Chiyuan Zhang, and Mengdi Wang. 2025a. Math-perturb: Benchmarking Ilms' math reasoning abilities against hard perturbations. *Preprint*, arXiv:2502.06453.
- Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and Ting Liu. 2025b. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *ACM Trans. Inf. Syst.*, 43(2).
- Takeshi Kojima, Shixiang (Shane) Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. In *Advances in Neural Information Processing Systems*, volume 35, pages 22199–22213. Curran Associates, Inc.
- Brenden Lake and Marco Baroni. 2018. Generalization without systematicity: On the compositional skills of sequence-to-sequence recurrent networks. In Proceedings of the 35th International Conference on Machine Learning, volume 80 of Proceedings of Machine Learning Research, pages 2873–2882. PMLR.
- Brenden M. Lake, Ruslan Salakhutdinov, and Joshua B. Tenenbaum. 2015. Human-level concept learning through probabilistic program induction. *Science*, 350(6266):1332–1338.
- Brenden M. Lake, Tomer D. Ullman, Joshua B. Tenenbaum, and Samuel J. Gershman. 2017. Building machines that learn and think like people. *Behavioral and Brain Sciences*, 40:e253.
- Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. 2022. Fantastically ordered prompts and where to find them: Overcoming fewshot prompt order sensitivity. In *Proceedings of the*

894

895

896

897

898

849

850

60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 8086–8098, Dublin, Ireland. Association for Computational Linguistics.

796

797

811

812

814

815

816

818

819

821

822

823

824

826

827

830

831

832

833

834

839

840

841

842

844

- Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. 2023. Self-refine: Iterative refinement with self-feedback. In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Suvir Mirchandani, Fei Xia, Pete Florence, brian ichter, Danny Driess, Montserrat Gonzalez Arenas, Kanishka Rao, Dorsa Sadigh, and Andy Zeng. 2023. Large language models as general pattern machines. In 7th Annual Conference on Robot Learning.
 - Iman Mirzadeh, Keivan Alizadeh, Hooman Shahrokhi, Oncel Tuzel, Samy Bengio, and Mehrdad Farajtabar. 2024. Gsm-symbolic: Understanding the limitations of mathematical reasoning in large language models. *Preprint*, arXiv:2410.05229.
 - Augustus Odena, Kensen Shi, David Bieber, Rishabh Singh, Charles Sutton, and Hanjun Dai. 2021. Bustle: Bottom-up program synthesis through learningguided exploration. *Preprint*, arXiv:2007.14381.
 - OpenAI. 2024a. Gpt-40 mini: advancing cost-efficient intelligence.
 - OpenAI. 2024b. Hello gpt-4o.
 - Linlu Qiu, Liwei Jiang, Ximing Lu, Melanie Sclar, Valentina Pyatkin, Chandra Bhagavatula, Bailin Wang, Yoon Kim, Yejin Choi, Nouha Dziri, and Xiang Ren. 2024. Phenomenal yet puzzling: Testing inductive reasoning capabilities of language models with hypothesis refinement. In *The Twelfth International Conference on Learning Representations*.
 - Joshua S Rule. 2020. *The child as hacker: Building more human-like models of learning*. Ph.D. thesis, MIT.
 - Mathias Sablé-Meyer, Kevin Ellis, Josh Tenenbaum, and Stanislas Dehaene. 2022. A language of thought for the mental representation of geometric shapes. *Cognitive Psychology*, 139:101527.
 - Joshua B. Tenenbaum, Charles Kemp, Thomas L. Griffiths, and Noah D. Goodman. 2011. How to grow a mind: Statistics, structure, and abstraction. *Science*, 331(6022):1279–1285.
- Lucas Tian, Kevin Ellis, Marta Kryven, and Josh Tenenbaum. 2020. Learning abstract structure for drawing by efficient motor program induction. In *Advances in Neural Information Processing Systems*, volume 33, pages 2686–2697. Curran Associates, Inc.

- Shangqing Tu, Chunyang Li, Jifan Yu, Xiaozhi Wang, Lei Hou, and Juanzi Li. 2024. Chatlog: Carefully evaluating the evolution of chatgpt across time. *Preprint*, arXiv:2304.14106.
- Ruocheng Wang, Eric Zelikman, Gabriel Poesia, Yewen Pu, Nick Haber, and Noah Goodman. 2024. Hypothesis search: Inductive reasoning with language models. In *The Twelfth International Conference on Learning Representations*.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023. Self-consistency improves chain of thought reasoning in language models. In *The Eleventh International Conference on Learning Representations*.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed H. Chi, Quoc V Le, and Denny Zhou. 2022. Chain of thought prompting elicits reasoning in large language models. In Advances in Neural Information Processing Systems.
- Zhaofeng Wu, Linlu Qiu, Alexis Ross, Ekin Akyürek, Boyuan Chen, Bailin Wang, Najoung Kim, Jacob Andreas, and Yoon Kim. 2024. Reasoning or reciting? exploring the capabilities and limitations of language models through counterfactual tasks. In Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), pages 1819–1862, Mexico City, Mexico. Association for Computational Linguistics.
- Zonglin Yang, Li Dong, Xinya Du, Hao Cheng, Erik Cambria, Xiaodong Liu, Jianfeng Gao, and Furu Wei. 2024. Language models as inductive reasoners. In Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers), pages 209–225, St. Julian's, Malta. Association for Computational Linguistics.
- Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu, Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang, Yulong Chen, Longyue Wang, Anh Tuan Luu, Wei Bi, Freda Shi, and Shuming Shi. 2023. Siren's song in the ai ocean: A survey on hallucination in large language models. *Preprint*, arXiv:2309.01219.
- Zhanke Zhou, Rong Tao, Jianing Zhu, Yiwen Luo, Zengmao Wang, and Bo Han. 2024. Can language models perform robust reasoning in chain-of-thought prompting with noisy rationales? In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.

900

903

904

905

906

907

908

909

911

912

913

914

915

916

917

918

919

924

928

929

930

932

934

936

937

938

941

943

945

946

Appendices A Details on Evaluation Pipeline

In this section, we provide detailed information on the evaluation pipeline, including the data construction and the performance assessment. For Arithmetic, we generate the base-7, base-8, and base-9 tasks by randomly sampling two two-digit numbers in the corresponding base, and then we check whether there is a carry-over in the addition process. If there is no carry-over, we regenerate the numbers. The noisy examples are generated in base-10. For Cryptography, we randomly select words of appropriate lengths from the NLTK Word Lists corpus¹, and then we encrypt the words using the Caesar, Atbash, and Keyboard ciphers. The noisy examples are generated by randomly replacing the letters in the output with other letters. For List Functions, we first write the corresponding rule functions for each task, and then we automatically generate the input data with appropriate lengths and ranges. The inputs are generated by randomly sampling numbers from a specific range with some constraints. . During the data synthesis process, we attach manual supervision to ensure that the generated data can correctly induce the rules. During the evaluation process, we use the inputs in the test set as the input for rule execution. We evaluate the model's performance using exact match. If the model's output is correct on all the test set examples, we consider the model to have successfully induced the rule. If the model fails to output a valid programmatic rule or the program contains an infinite loop or errors, we consider it a failure.

B Experimental Details

B.1 Experimental Settings

For robustness under different noise levels, we run each experiment three times and report the mean and standard deviation of the results to avoid randomness. Except for the self-consistency (SC) and sample-steered rule refinement (SRR) that require diverse generations, we set the temperature to 0.0 for all models to ensure reproducibility. For SC and SRR, the temperature is set to 0.7, consistent with the original work of SC (Wang et al., 2023). The positions of noise in seen examples are random to avoid positional bias (Lu et al., 2022). In the implementation, we choose 2 subsets for SRR by splitting the seen examples into two parts. The number of iterations is set to 3 for both SR and SRR.

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

972

All experiments are conducted through the official OpenAI² and DeepSeek³ API platform. For the GPT-4o-mini and GPT-4o models, we spend about 200 USD in total, and for the DeepSeek-V3 and DeepSeek-R1 models, we spend about 50 USD in total.

B.2 Prompts and Failure Cases

For the Direct Output setting, we restrict the model to output the rule directly without any additional output, as shown in Table 7. For the chain-ofthought and self-consistency setting, we use the instruction in Kojima et al. (2022). For the selfrefine and sample-steered rule refinement setting, we use the chain-of-thought prompt as the initial prompt; the iterative prompts are shown in Table 9 and Table 10, respectively.

We provide some representative failure cases in the evaluation. For the Arithmetic task, the model fails to solve the base-7 and base-9 tasks and misinterprets the rule as the base-8 addition or the decimal sum with a constant. For the Cryptography task, the model fails to solve the Atbash and tries to explain it as shifts. The responses of the models are shown in Table 11.

²https://openai.com/api/

³https://platform.deepseek.com/usage

Prompt for Direct Output

Please generate a rule that maps the following inputs to their corresponding outputs using a Python function. The input is a list of integers. The output is also a list of integers. Note that some examples may be wrong, and you should take this into account when proposing the rule. {examples} Please format your Python function as follows: "'python def fn(x): # Your code here "'

Your response should only include the function definition, not the function call or any other information.

Table 7: The prompt for the Direct Output setting, exemplified by the List Functions task. The {examples} in the prompt is replaced by input-output pairs when conducting the experiments.

Prompt for Chain-of-Thought

Please generate a rule that maps the following inputs to their corresponding outputs using a Python function. The input is a list of integers. The output is also a list of integers. Note that some examples may be wrong, and you should take this into account when proposing the rule. {examples} Please format your Python function as follows: ""python

def fn(x):

Your code here

Think step-by-step and explain your reasoning. Your response should include your thought process and the function definition without the function call.

Table 8: The prompt for the Chain-of-Thought setting, exemplified by the List Functions task.

Prompt for Feedback Generation in Self-Refine

You have generated a rule that maps the following inputs to their corresponding outputs using a Python function. The input is a list of integers. The output is also a list of integers. {examples} In the last step, your rule is: "'python {rule} ""

Give some feedback on the rule you have generated, like how can it be improved, what is wrong with it, etc. Your response should only include the feedback. If you think the rule is good enough, your response should be "NO FEEDBACK" without other information. Note that some examples may be wrong, and you should take this into account when proposing the feedback.

Iteration

You have generated a rule that maps the following inputs to their corresponding outputs using a Python function. The input is a list of integers. The output is also a list of integers. Note that some examples may be wrong, and you should take this into account when proposing the rule. {examples} In the last step, your rule is: "python {rule} The feedback you have given is: {feedback} Generate a new rule that maps the given inputs to their corresponding outputs using a Python function. Please format your rule as follows: "python def fn(x): # Your code here Think step-by-step and explain your reasoning. Your response should include your thought process and the function definition without the function call.

Table 9: The prompt for the Self-refine setting, exemplified by the List Functions task. The {examples}, {rule} and {feedback} in the prompt are replaced by input-output pairs, rule, and feedback when conducting the experiments.

Prompt for Iteration in SRR

You have generated a rule that maps the following inputs to their corresponding outputs using a Python function. The input is a list of integers. The output is also a list of integers. Note that some examples may be noisy, and you should take this into account when proposing the rule. In the last step, your rule is "'python {rule} "' But this rule is not correct. It works for the following examples: {right_examples} However, it does not work for the following examples: {wrong_examples} Generate a new rule that maps the given inputs to their corresponding outputs using a Python function. Please format your rule as follows: "'python def fn(x): # Your code here "' Think step-by-step and explain your reasoning. Your response should include your thought process and the function

I hink step-by-step and explain your reasoning. Your response should include your thought process and the function definition without the function call. You can either modify the existing rule or propose a new one.

Table 10: The prompt for the Sample-steered Rule Refinement setting, exemplified by the List Functions task. The {rule}, {right_examples}, and {wrong_examples} in the prompt are replaced by the rule, right examples, and wrong examples when conducting the experiments.

Failure Cases for Arithmetic

Octonary Confusion Error (in base-7):

To solve this problem, we need to identify a pattern or rule that maps the given input expressions to their corresponding outputs. Let's analyze the examples provided: First, let's verify if the outputs are simply the sum of the inputs in base 10:- For '34 + 62', the sum is '96', but the output is '126'.- For '56 + 30', the sum is '86', but the output is '116'.- For '22 + 25', the sum is '47', but the output is '50'... Clearly, the outputs are not the direct sums of the inputs in base 10. Let's consider the possibility that the numbers are in base 8 (octal):- Convert the inputs to base 10, perform the addition, and then convert the result back to base 8. Let's verify this hypothesis: 1. 34 + 62 in base 8: - '34' in base 8 is '3*8 + 4 = 28' in base 10. - '62' in base 8 is '6*8 + 2 = 50' in base 10. - Sum in base 10: '28 + 50 = 78'. - Convert '78' to base 8: '78 // 8 = 9' remainder '6', so '78' in base 8 is '116'.2. 56 + 30 in base 8: - '56' in base 8 is '5*8 + 6 = 46' in base 10. - '30' in base 8 is '3*8 + 0 = 24' in base 10. - Sum in base 10: '46 + 24 = 70'. - Convert '70' to base 8: '70 // 8 = 8' remainder '6', so '70' in base 8 is '106'. 3. 22 + 25 in base 8: - '22' in base 8 is '2*8 + 2 = 18' in base 10. - '25' in base 8 is '2*8 + 5 = 21' in base 10. - Sum in base 10: '18 + 21 = 39'. - Convert '39' to base 8: '39 // 8 = 4' remainder '7', so '39' in base 8 is '47'. Following this pattern, we can see that the outputs are indeed the sums of the inputs when interpreted as base 8 numbers and then converted back to base 8. Let's implement this in a Python function...

Decimal Confusion Error (in base-9):

To solve this problem, we need to identify a consistent rule or pattern that maps the given input expressions to their corresponding outputs. Let's analyze the provided examples: 1.68 + 68 = 147: The correct sum is 136, but the output is 147. 2.80 + 42 = 132: The correct sum is 122, but the output is 132. 3.76 + 17 = 104: The correct sum is 93, but the output is 104. 4.38 + 14 = 53: The correct sum is 52, but the output is 53. 5.43 + 84 = 137: The correct sum is 127, but the output is 137. 6.48 + 55 = 114: The correct sum is 103, but the output is 114. 7.72 + 56 = 138: The correct sum is 128, but the output is 138. 8.34 + 46 = 81: The correct sum is 80, but the output is 81. 9.88 + 15 = 114: The correct sum is 103, but the output is 117. From these examples, it appears that the output is consistently 10 more than the correct sum of the two numbers. This suggests that the rule is to calculate the sum of the two numbers and then add 10 to the result. Let's define a Python function to implement this rule...

Failure Case for Cryptography

Shifts Error:

To solve this problem, we need to identify a consistent rule or transformation that maps each input string to its corresponding output string. Let's analyze the given examples to find a pattern: 1. cuissard -> xfrhhziw 2. tripsis -> girkhrh 3. Maccabaeus -> Nzxzyzvfh 4. jocundness -> qlxfmwmvhh 5. chiro -> xsril 6. coapt -> xlzkg 7. Phecda -> Ksvxwz 8. unifilar -> fmrurozi 9. Cordaites -> Xliwzrgvh 10. Wendell -> Dvmwvoo ### Observations: 1. Character Shifts : The transformation seems to involve shifting characters. For example, 'c' in "cuissard" becomes 'x' in "xfrhhziw". This suggests a possible Caesar cipher or similar character shift. 2. Case Sensitivity : The case of the letters is preserved. For instance, 'M' in "Maccabaeus" becomes 'N' in "Nzxxzyzvfh", and 'P' in "Phecda" becomes 'K' in "Ksvxwz". 3. Consistent Shifts : By examining the shifts: - 'c' to 'x' is a shift of 23 positions backward in the alphabet. - 't' to 'g' is a shift of 13 positions backward. - 'M' to 'N' is a shift of 1 position forward. - 'j' to 'q' is a shift of 7 positions forward. - 'P' to 'K' is a shift of 21 positions backward. - 'u' to 'f' is a shift of 11 positions backward. - 'C' to 'X' is a shift of 21 positions backward. - 'W' to 'D' is a shift of 21 positions backward. 4. Pattern Recognition : The shifts are not consistent across all examples, suggesting that some examples might be incorrect or that multiple rules are being applied. However, a common shift that appears frequently is 21 positions backward (e.g., 'c' to 'x', 'P' to 'K', 'C' to 'X', 'W' to 'D'). ### Hypothesis: Given the frequent appearance of a 21-position backward shift, let's hypothesize that the rule involves shifting each character 21 positions backward in the alphabet. This would mean: - For lowercase letters, 'a' becomes 'v', 'b' becomes 'w', ..., 'z' becomes 'u'. - For uppercase letters, 'A' becomes 'V', 'B' becomes 'W', ..., 'Z' becomes 'U'. ### Implementation: Let's implement this hypothesis in a Python function...

Table 11: The failure cases of the models on the Arithmetic and Cryptography tasks.